

ADAPTIVE READER-WRITER LOCK

CROSS REFERENCE TO RELATED APPLICATION(S)

The present application is a continuation of U.S. Patent Application Serial No.
5 09/741,679, now pending.

BACKGROUND OF THE INVENTION

Technical Field

10 This invention relates to resource locking in computer systems and more specifically to a method of dynamically selecting an optimal lock mode. Both units within a central processing system and system-wide measurements are maintained, and based upon these measures an optimal locking mode is determined.

Description Of The Prior Art

15 Multiprocessor systems contain multiple processors (also referred to herein as CPUs) that can execute multiple processes or multiple threads within a single process simultaneously in a manner known as parallel computing. In general, multiprocessor systems execute multiple processes or threads faster than conventional single processor systems, such as personal computer, that execute programs sequentially. The actual
20 performance advantage is a function of a number of factors, including the degree to which parts of a multithreaded process and/or multiple distinct processes can be executed in parallel and the architecture of the particular multiprocessor system. The degree to which processes can be executed in parallel depends, in part, on the extent to which they compete for exclusive access to shared memory resources.

Shared memory multiprocessor systems offer a common physical memory address space that all processors can access. Multiple processes therein, or multiple threads within a process, can communicate through shared variables in memory which allow the processes to read or write to the same memory location in the computer system. Message passing multiprocessor systems, in contrast to shared memory systems, have a separate memory space for each processor. They require processes to communicate through explicit messages to each other.

The architecture of shared memory multiprocessor systems may be classified by how memory is physically organized. In distributed shared memory (DSM) machines, the memory is divided into modules physically placed near one or more processors, typically on a processor node. Although all of the memory modules are globally accessible, a processor can access local memory on its node faster than remote memory on other nodes. Because the memory access time differs based on memory locations, such systems are also called non-uniform memory access (NUMA) machines. In centralized shared memory machines, the memory is physically in one location. Centralized shared memory computers are called uniform memory access (UMA) machines because the memory is equidistant in time from each of the processors. Both forms of memory organization typically use high-speed cache in conjunction with main memory to reduce execution time.

The use of NUMA architecture to increase performance is not restricted to NUMA machines. A subset of processors in a UMA machine may share a cache. In such an arrangement, even though the memory is equidistant from all processors, data can circulate among the cache-sharing processors faster (*i.e.*, with lower latency) than among the other processors in the machine. Algorithms that enhance the performance of NUMA machines can be applied to any multiprocessor system that has a subset of processors with lower latencies. These include not only the noted NUMA and shared cache machines, but also machines where multiple processors share a set of bus-interface logic as well as machines with interconnects that “fan out” (typically in hierarchical fashion) to the processors.

A significant issue in the design of multiprocessor systems is process synchronization. The degree to which processes can be executed in parallel depends in part on the extent to which they compete for exclusive access to shared memory resources. For example, if two processes A and B are executing in parallel, process B might have to wait for process A to write a value to a buffer before process B can access it. Otherwise, a race condition could occur, where process B might access the buffer while process A was part way through updating the buffer. To avoid conflicts, process synchronization mechanisms are provided to control the order of process execution. These mechanisms include mutual exclusion locks, condition variables, counting semaphores, and reader-writer locks. A mutual exclusion lock allows only the processor holding the lock to execute an associated action. When a processor requests a mutual exclusion lock, it is granted to that processor exclusively. Other processors desiring the lock must wait until the processor with the lock releases it. To address the buffer scenario described above, both processes would request the mutual exclusion lock before executing further. Whichever process first acquires the lock then updates (in the case of process A) or accesses (in the case of process B) the buffer. The other processor must wait until the first processor finishes and releases the lock. In this way, the lock guarantees that process B sees consistent information, even if processors running in parallel execute processes A and B.

For processes to be synchronized, instructions requiring exclusive access can be grouped into a critical section and associated with a lock. When a process is executing instructions in its critical section, a mutual exclusion lock guarantees no other processes are executing the same instructions. This is important where processes are attempting to change data. However, such a lock has the drawback in that it prohibits multiple processes from simultaneously executing instructions that only allow the processes to read data. A reader-writer lock, in contrast, allows multiple reading processes (“readers”) to access simultaneously a shared resource such as a database, while a writing process (“writer”) must have exclusive access to the database before performing any updates for consistency. A practical example of a situation appropriate for a reader-writer lock is a TCP/IP routing structure with many readers and an occasional update of

the information. Recent implementations of reader-writer locks are described by Mellor-Crummey and Scott (MCS) in “Scalable Reader-Writer Synchronization for Shared-Memory Multiprocessors,” *Proceedings of the Third ACM SIGPLAN Symposium on Principles and Practices of Parallel Programming*, pages 106-113 (1991) and by Hsieh and Weihs in “Scalable Reader-Writer Locks for Parallel Systems,” *Technical Report MIT/LCS/TR-521* (November 1991).

The basic mechanics and structure of reader-writer locks are well known. In a typical lock, multiple readers may acquire the lock, but only if there are no active writers. Conversely, a writer may acquire the lock only if there are no active readers or another writer. When a reader releases the lock, it takes no action unless it is the last active reader; if so, it grants the lock to the next waiting writer. A centralized reader-writer lock mode is a formative of a reader-writer lock that uses a single data structure to control access to the lock. There are many implementations of this locking primitive. One of the simplest formative uses a set of counters guarded by a simple test-and-set spinlock. The counters count the number of readers holding the lock, the number of readers waiting for access to the lock, the number of writers holding the lock (which must be either one or zero), and the number of writers waiting on the lock. The readers and writers go through a decision process based upon the counter values. Accordingly, this mode is optimal for high update rates wherein read side critical sections are lengthy. The simple test and set lock is a form of an exclusive lock. Other types of exclusive locks include test and set; test and test and set; queued lock; ticket lock; and quad aware lock.

A drawback of prior reader-writer locks is undesired memory contention, whereby multiple processors modify a single data structure in quick succession, possibly while other processors are spinning on said single data structure. The resulting cache misses can severely degrade performance. The drawback has been partially addressed in more recent locking schemes such as the ones described by Hsieh and Weihs. Their static locking algorithm allocates one semaphore per processor, stored in memory local to the processor. An additional semaphore acts as a gate on the writers. To acquire a static lock, a reader need only acquire its local semaphore, greatly reducing the amount of spinning. However, a writer must still acquire all of the semaphores of which there is

now one for each processor, and the additional semaphore. When releasing a static lock, a reader simply releases its local semaphore, and a writer releases all of the semaphores. The lock thus offers an improvement over prior locks in that the readers do not interfere with each other and readers do not have to go over the system interconnect to acquire a lock. However, the fact that readers never interfere means that writers must do a substantial amount of work in systems with many processors. When even a few percent of the requests are writes, the throughput suffers dramatically because a writer must acquire a semaphore for every processor on every node to successfully acquire the lock. To overcome this problem, their dynamic locking scheme attempts to reduce the number of semaphores a writer must acquire by keeping track of active readers in a single memory location and acquiring only semaphores associated with these readers. The scheme uses a variety of mutual exclusion locks to accomplish this.

“Reactive Synchronization Algorithms for Multiprocessor” by Beng-Hong Lim et al. describes an adaptive exclusive lock which teaches the performance benefits of selecting synchronization protocols in response to the level of contention. The disclosure teaches switching between a simple test-and-set spinlock and a queued lock, both of which are exclusive locks. Beng-Hong Lim et al. teach dynamically switching to simple test and set spinlock at low contention and to queued lock at high contention, thereby using each of these locking modes when it operates most effectively. However, Beng-Hong Lim et al. does not teach a method of dynamically selecting a locking mode wherein different modes may be beneficial for differing ratios of read and write requests.

Accordingly, there is a need for a computer system comprising multiple processors, means for determining an optimal locking mode, and means for switching among locking modes based upon the determination.

25

SUMMARY OF THE INVENTION

It is therefore an object of the invention to provide a method of dynamically selecting a lock mode in a multiprocessor computer system. It is a further object of the

invention to switch among locking modes based upon system parameters and measurements to provide an efficient operating mode.

A first aspect of the invention is a method of dynamically determining a lock mode in a multiprocessor computer system. First and second system-wide measures of read and write acquisitions are maintained. A lock mode manager is implemented to select an optimal lock mode based upon at least some of these measures. Optionally, the system is switched to that lock mode from another lock mode. The lock mode is preferably, but not necessarily, selected from the group consist of: a distributed reader-writer lock mode, a centralized reader-writer lock mode, and an exclusive lock mode.

A second aspect of the invention is a multiprocessor computer system including first and second system-wide measures of read and write acquisitions. A lock mode manager is implemented to select an optimal locking mode, responsive to at least some of those measures. The lock mode is preferably, but not necessarily, selected from the group consist of: a distributed reader-writer lock mode, a centralized reader-writer lock mode, and an exclusive lock mode.

A third aspect of the invention is an article in a multiprocessor system is provided comprising a computer-readable signal bearing medium. The article includes means for maintaining first and second system-wide measures of read and write acquisitions. In addition, means in the medium are provided for selecting a lock mode based upon some of the measures. The lock mode is preferably, but not necessarily, selected from the group consist of: a distributed reader-writer lock mode, a centralized reader-writer lock mode, and an exclusive lock mode.

Other features and advantages of this invention will become apparent from the following detailed description of the presently preferred embodiment of the invention, taken in conjunction with the accompanying drawings.

BRIEF DESCRIPTION OF THE DRAWINGS

FIG. 1 is a flow chart illustrating selection of an optimal locking mode according to the preferred embodiment of this invention, and is suggested for printing on the first page of the issued patent; and

FIG. 2 is a data structure diagram of the system measurements.

5

DESCRIPTION OF THE PREFERRED EMBODIMENT

Overview

Technical Background

Maintaining low lock contention is essential to attaining high performance in parallel programs. However, even programs with negligible lock contention can suffer severe performance degradation due to memory latencies incurred when accessing shared data that is frequently modified. This is due to the high cost of memory latency compared to instruction execution overheads. Memory latencies are incurred for shared data structures in addition to the locks themselves. Accordingly, it is desirable and efficient to operate in different locking modes under differing levels of operation.

It has been ascertained that there are essentially three desirable forms of locking modes or primitives that provide efficient modes of operation for differing levels of contention, two formations of a reader-writer lock mode, and one form of an exclusive lock. Although there are many locking primitives which have been developed and maintained, these three primitives provide efficient forms of locks that address issues with high performance in parallel programs. In general, a reader-writer lock is a lock that has a reader mode and a writer mode. Multiple readers can hold the lock simultaneously, but only a single writer is allowed to hold the lock at a given time. As such, readers exclude writers, writers exclude readers, and writers exclude writers. The two forms of reader-writer locks employed within the dynamic lock primitive of the preferred embodiment are a distributed reader-writer spinlock and a centralized reader-writer lock. In addition, a simple test-and-set spinlock is incorporated as an alternative locking primitive. Although the test-and-set spinlock is not a reader-writer lock, it may typically

be used where one would use a reader-writer lock, giving up the ability for readers to execute in parallel. Accordingly, although the dynamic lock primitive is taught herein with specific lock modes, alternative lock modes may be employed where and when desirable.

5 A distributed reader-writer spinlock is constructed by maintaining a separate simple spinlock per CPU, and an additional simple spinlock to serialize write-side accesses. Each of these locks is in its own cache line in order to prevent false sharing. To read-acquire the distributed reader-writer spinlock, the readers acquire the lock corresponding to their CPU. To read release a distributed reader-writer spinlock, a CPU releases its lock. To write acquire the distributed reader-writer spinlock, a CPU first acquires the writer gate, then each of the CPU's spinlocks. To write-release the lock, the per-CPU locks are released, and then the writer gate is released. The distributed reader-writer lock performs well when the ratio of write acquisitions to total acquisitions is low. However it is a slower write-side lock mode than the simple test-and-set spinlock.

10 Accordingly, a distributed reader-writer spinlock is optimal for a high quantity of read accesses with short hold times, and should be used only if less than 1 out of the number of CPUs of the accesses are write-side accesses.

15 A centralized reader-writer lock mode is a formative of a reader-writer lock that uses a single data structure to control access to the lock. There are many implementations of this locking primitive. One of the simplest formative uses a set of counters guarded by a simple test-and-set spinlock. The counters count the number of readers holding the lock, the number of readers waiting for access to the lock, the number of writers holding the lock (which must be either one or zero), and the number of writers waiting on the lock. The readers and writers go through a decision process based upon 20 the counter values. Accordingly, this mode is optimal for high update rates wherein read side critical sections are lengthy.

25 A simple test and set primitive allows a process to request a lock by repeatedly executing instructions on a boolean flag until it successfully changes the flag from false to true. Simple test-and-set spin locks have been used for decades and are often the primitive of choice for low contention due to their extremely short code-path lengths.

However, they perform poorly under high read-side contention, because they do not allow read-side acquisitions to proceed in parallel. Accordingly, the simple test-and-set mode is optimal when there are frequent write-side accesses.

Most operating systems are preprogrammed to operate in a singular reader-writer lock mode. Dynamically choosing operating modes presents a challenge in that a reactive algorithm needs to select and change modes efficiently and has to allow for the possibility that multiple processes may be executing different protocols at the same time. Frequently, the best operating mode depends on the level of contention, as well as the magnitude of the fraction of total acquisitions that are write acquisitions. In designing reactive algorithms that dynamically select among multiple synchronization protocols several issues must be addressed, including: how the algorithm efficiently detects which modes to use, how the algorithm changes modes correctly and efficiently, and when the algorithm should change protocols. Accordingly, in providing efficient mechanisms for a dynamic protocol selection that multiple processes may be trying to execute, the main considerations are execution of the synchronization operation and maintaining constant agreement as to which protocol to use without incurring significant overhead.

Each of the implementations of reader-writer locks and the exclusive test-and-set spinlock described herein and known in the art are optimal under different operation conditions. For example, a distributed reader-writer spinlock works best under highly read-intensive workloads, a centralized reader-writer spinlock works best under moderately read-intensive workloads with large read-side critical sections, and a simple test-and-set spinlock works best under write-intensive workloads. It is known in the art, that the programmer using a locking mode must select the locking mode to be implemented when coding. Accordingly, if the workload shifts to different operating conditions, the single locking algorithm pre-selected can become non-optimal.

Fig. 1 is a flow chart 10 illustrating the path followed for dynamically selecting a reader-writer lock mode in a multiprocessor. As described above, there are three general forms of a reader-writer lock mode, a distributed reader-writer lock mode 12, a centralized reader-writer lock mode 14, and a simple test-and-set lock mode 16. Each of these operating modes are beneficial and more efficient under different operating

conditions. The essence of the method herein is to ascertain which locking mode is a more efficient mode under different operating conditions and switching to that optimal mode under specific operating conditions.

There are different measurements acquired to determine when it is desirable to switch among operating modes. Fig. 2 is a diagram 20 illustrating the different measurements that are acquired and where they are stored. Each CPU 24 in the system maintains a measure of the quantity of read-side acquisitions 26 of that CPU. In addition, each CPU maintains a measure of the read-hold duration 28 for that CPU. The measure 24 and 26 may be in the form of a counter, or an alternative measurement for maintaining the quantity of read-side acquisitions on each CPU. In an alternative embodiment, information is collected by a unit within a computer processing system, wherein the unit may be in the form of a thread, a process, a transaction, a co-routine, a thread in a multi-threaded architecture, and a task, and the counters measure read-side and read-hold duration by the appropriate unit. The read-hold duration 28 may be in the form of a digital filter, or an alternative mathematical algorithm. The digital filter may compute a weighted average, a sliding window average, a finite impulse response, and alternative computational operations that adjust the ability to respond to current events as opposed to historic events of the operating mode. External to the per-CPU measurements is a central data structure 30 which operates as a central location for storing system wide measurements and maintaining a system wide accounting. The central data structure maintains information for a measurement of the quantity of read-side acquisitions in the system 32 and the quantity of write-side acquisitions in the system 34. In addition, the central data structure 30 maintains a measurement of the read and write acquisitions 36 and 38, respectively, and a measurement of read-hold durations 40. Accordingly, the measurements maintained in each CPU, 26 and 28, and the measurements maintained in the central data structure 32-40 are used to compute and determine an optimal operating mode for the operating system and to enhance efficiency for switching among operating modes.

The central measure of read acquisitions 36 may be in the form of a digital filter, or an alternative mathematical algorithm, that contributes to the accuracy of calculating

the optimal operating mode. The measure 36 in the form of a digital filter is a computational means for determining recent events in the quantity of read-hold acquisitions in the system. The digital filter may compute a weighted average, a sliding window average, a finite impulse response, and alternative computational operations that 5 adjust the ability to respond to current events as opposed to historic events of the operating mode.

The central measure of write acquisitions 38 may also be in the form of a digital filter, or an alternative mathematical algorithm, that contributes to the accuracy of calculating the optimal operating mode. The measure 38 in the form of a digital filter is a 10 computational means for determining recent events in the quantity of write-hold acquisitions in the system. The digital filter may compute a weighted average, a sliding window average, a finite impulse response, and alternative computational operations that adjust the ability to respond to current events as opposed to historic events of the operating mode. The central data structure does not store information pertaining to the 15 write-hold per CPU. The write-hold requires a manipulation of the central data structure regardless of the operating mode. Therefore, it is only necessary to store such information in the central data structure.

Finally, the central data structure stores and maintains a measurement of the read-hold duration 40. This measure computes a time interval for which the CPUs in the 20 system are using read-holds. The measure 40 may come in the form of a digital filter, which may be a weighted average, a sliding window average, a finite impulse response, or an alternative computations operation that adjust the ability to determine the centralized reader-write lock mode 14 and the simple test-and-set lock mode 16 for acquiring a write lock. Accordingly, the read-hold duration measure is used to determine 25 an optimal and efficient time for switching between the two operating modes that are desirable for write-hold modes.

When a CPU in the system write-acquires the lock, the CPU measure of read-side acquisitions 26 and the quantity of read-side acquisitions are applied to the measure 32 in the central data structure and the CPU measure 26 is reset to zero. A measure of write-side acquisitions 34 is maintained by a discrete unit of information, which is periodically 30

updated. In a preferred embodiment, the update method may be a computer implemented algorithm or computer code. At such time as the measure 32 of read-side acquisitions in the central data structure 30 is sufficiently greater than the measure 34 of write-side acquisitions, and if the lock mode is operating outside of the distributed reader-writer lock mode 12, then the lock mode switches into the distributed reader-writer lock mode 12. The transfer into this mode 12 is illustrated in Fig. 1. As shown in Fig. 1, both the centralized reader-writer lock mode 14 and the simple test-and-set lock mode 16 are desirable for write intensive workloads. However, at such time as the measure of read acquisitions exceeds the measure of write acquisitions, it may be desirable to switch to a distributed reader-write lock mode for more efficient operation. If the system were already operating in a distributed reader-writer lock mode, then it would remain in this mode until such time as it is determined optimal to switch to an alternative lock mode to enhance operating efficiency. However, regardless of which locking mode the system is operating under and which mode it may be optimal to switch to, a lock must be held on the central lock state before a switching of lock modes may occur. This prevents two CPUs from attempting to switch states at the same time. The lock on the central lock is held during any acquisition and release for the simple test-and-set spinlock and the centralized read-writer lock modes, and for write acquisition and release for the distributed reader-writer lock mode.

Having defined all of the measurement maintained by each CPU and the central data structure, these measurements may now be mathematically manipulated to determine optimal times to switching operating modes. The following terms define the measurements and mathematics employed: D_r is the measure of duration of read acquisitions, n_r is the measure of read acquisitions, n_c is the number of CPUs in the system, n_w is the measure of the number of write acquisitions, d_r is D_r/n_r , f is $n_w / (n_r + n_w)$, t_f is the threshold for f shown as $(1/(2*n_c))$, t_D is the threshold for the accumulated duration of read acquisitions shown as one half of a CPU over a period of time, L_m is the amount of time it takes for a memory request to be satisfied given that memory requested was recently modified by another CPU, i.e. memory latency, and t_d is the threshold for d_r , defined as $L_m * 5$.

If the system lock mode is operating in a distributed reader-writer lock mode 12, it is presumed that the system is operating under a significant quantity of read-side acquisitions, and will switch from this mode to an alternate mode when servicing a write-side access. At such time as $f > t_f$ then it may be optimal to switch to either the 5 centralized reader-writer lock mode 14 or the simple test-and-set spinlock mode 16. The transfer into the centralized reader-writer lock mode 14 or the simple test-and-set spinlock mode 16 is illustrated in Fig. 1. As shown in Fig. 1, both the centralized reader-writer lock mode 14 and the simple test-and-set lock mode 16 are desirable for write intensive workloads, and the distributed reader-writer lock mode is desirable for 10 read intensive workloads. Accordingly, the preferred embodiment considers switching from a distributed reader-writer lock mode 12 only when servicing a write-side access.

If the system is operating in a centralized reader-writer lock mode 14, it may be desirable to switch to another operating mode under certain conditions and circumstances. For example, if $f < t_f$ it may be desirable to switch to a distributed reader-writer lock mode 12. Alternatively, if $d_r < t_d$ or if $D_r < t_D$ then it may be desirable to switch to a simple test-and-set spinlock. Finally, if the system is operating under the simple test-and-set spinlock mode, it may be desirable to switch to another operating mode under certain conditions and circumstances. For example, if $f < t_f$ then it may be 15 desirable to switch to a distributed reader-writer lock mode. However, if $d_r > t_d$ and $D_r > t_D$, then it may be desirable to switch to a centralized reader-writer lock mode. Accordingly, the process of switching operation modes is based upon mathematical evaluations from the values stored and defined in each CPU as well as the central data 20 structure, as shown in Fig. 2.

In an alternative embodiment, the system may only have the option of switching 25 between a distributed reader-writer lock mode and a centralized lock, wherein the centralized lock may be a simple test-and-set spinlock or a centralized reader-writer lock. If the system is operating in a distributed reader-writer lock, write-side access has been requested, and $f > t_f$, then it may be desirable to switch to a centralized lock mode. However, if the system is operating in a centralized lock mode and $f < t_f$, then it may be 30 desirable to switch to a distributed reader-writer lock mode. When operating between the

two modes noted herein, the mathematical formulae are less complex. However, the price of simplicity is less-optimal performance under some workloads.

Fig. 1 is a flow chart illustrating how the system switches among operating modes when it is determined that switching would provide more optimal operating conditions.

As shown in Fig. 1, if the system is operating in a distributed reader-writer lock mode 12, it may be desirable to switch to another operating mode when write-side acquisitions exceed a minimum threshold. For example, if the measure of write acquisitions is greater than one from a total of write acquisitions, then the desirable operating mode is the centralized reader-writer lock mode 14. However, if the system is operating in a centralized reader-writer lock mode 14, it may be desirable to switch to a simple test-and-set lock mode 16 when the read-hold duration measure is below a defined threshold. Alternatively, if the system is operating in a centralized reader-writer lock mode 14 and the read-hold measure exceeds a defined threshold, it is desirable to switch to a distributed reader-writer lock mode 12. Accordingly, the dynamic locking primitive utilizes predefined parameter for determining optimal switching opportunities.

Advantages Over The Prior Art

It is known in the art to implement an adaptive exclusive lock for an operating system. In the prior art, the operating mode may switch between a simple test-and-set spinlock and a queued lock. Both of these are exclusive locking modes and do not address a dynamic lock primitive for reader-writer spinlocks. However, the method of dynamically switching among lock modes for an operating system allows the system to switch under optimal conditions to provide an efficient use of time and memory. There is overhead involved with conducting a switch, and as such switching among nodes should only be induced when it has been determined to be optimal to provide a more efficient operation.

Alternative Embodiments

It will be appreciated that, although specific embodiments of the invention have been described herein for purposes of illustration, various modifications may be made without departing from the spirit and scope of the invention. In particular, the system maintains and calculates secondary measurements for use in determining optimal modes
5 of operation and an optimal opportunity for switching to an alternative locking mode. The secondary measurements come in alternative forms depending upon the needs ascertained. For example, it may be desirable to utilize a digital filter for such measurements. Alternatively, a weighted average or another statistical format may be desirable for determining accurate secondary measurements. In addition, the method
10 may be applied to a computer system and article of manufacture for use with the method of dynamically selecting among locking modes. Some systems and workloads may require separate thresholds be used to determine when to switch from versus to a given lock mode in order to provide hysteresis. Accordingly, the scope of protection of this invention is limited only by the following claims and their equivalents.